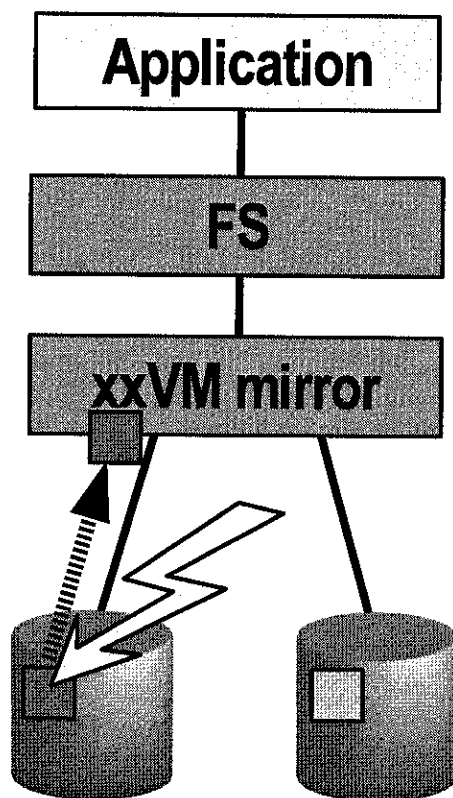
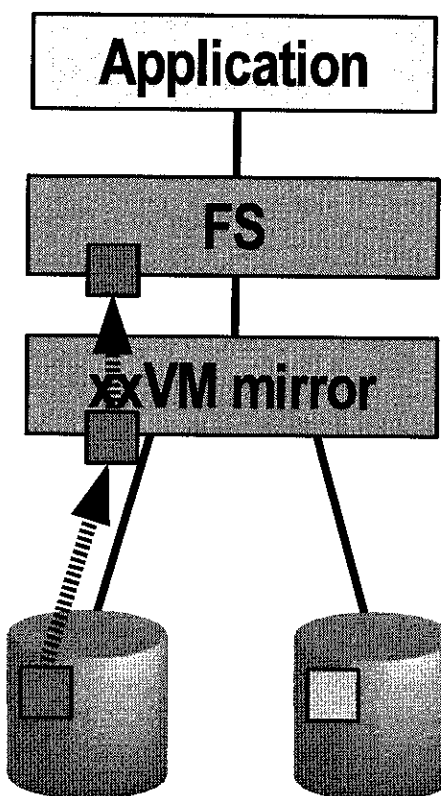


Traditional Mirroring

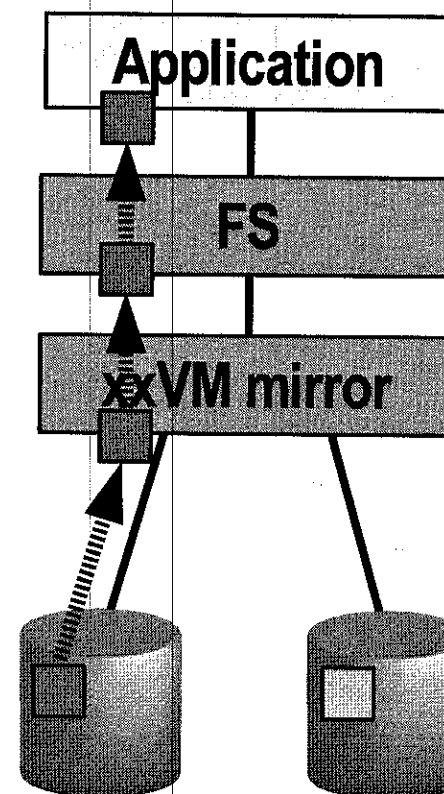
1. Application issues a read.
Mirror reads the first disk, which has a corrupt block. It can't tell.



2. Volume manager passes bad block up to filesystem.
If it's a metadata block, the filesystem panics. If not...

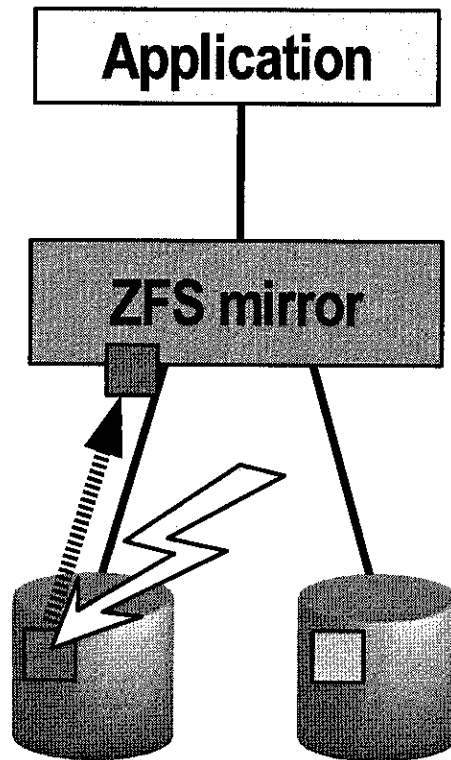


3. Filesystem returns bad data to the application.

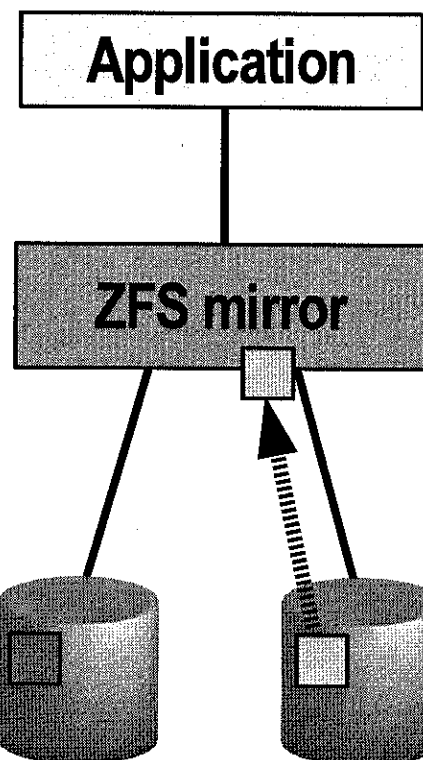


Self-Healing Data in ZFS

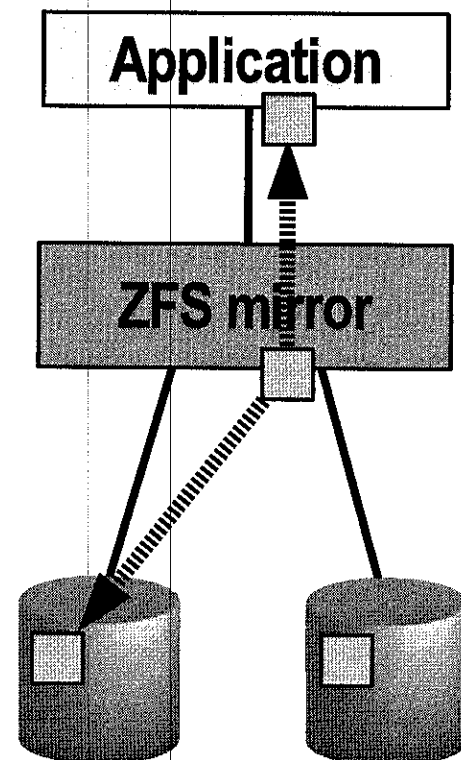
1. Application issues a read. ZFS mirror tries the first disk. Checksum reveals that the block is corrupt on disk.



2. ZFS tries the second disk. Checksum indicates that the block is good.



3. ZFS returns good data to the application and repairs the damaged block.



Traditional RAID-4 and RAID-5

- Several data disks plus one parity disk

$$\text{disk} \wedge \text{disk} \wedge \text{disk} \wedge \text{disk} \wedge \text{disk} = 0$$

- Fatal flaw: partial stripe writes
 - Parity update requires read-modify-write (slow)
 - Read old data and old parity (two synchronous disk reads)
 - Compute new parity = new data \wedge old data \wedge old parity
 - Write new data and new parity
 - Suffers from *write hole*: $\text{disk} \wedge \text{disk} \wedge \text{disk} \wedge \text{disk} \wedge \text{disk} = \text{garbage}$
 - Loss of power between data and parity writes will corrupt data
 - Workaround: \$\$\$ NVRAM in hardware (i.e., don't lose power!)
- Can't detect or correct silent data corruption

RAID-Z

- **Dynamic stripe width**
 - Variable block size: 512 – 128K
 - Each logical block is its own stripe
- **Both single- and double-parity**
- **All writes are full-stripe writes**
 - Eliminates read-modify-write (it's fast)
 - Eliminates the RAID-5 write hole (no need for NVRAM)
- **Detects and corrects silent data corruption**
 - Checksum-driven combinatorial reconstruction
- **No special hardware – ZFS loves cheap disks**

Disk		A	B	C	D	E
LBA						
0		P ₀	D ₀	D ₂	D ₄	D ₆
1		P ₁	D ₁	D ₃	D ₅	D ₇
2		P ₀	D ₀	D ₁	D ₂	P ₀
3		D ₀	D ₁	D ₂	P ₀	D ₀
4		P ₀	D ₀	D ₄	D ₈	D ₁₁
5		P ₁	D ₁	D ₅	D ₉	D ₁₂
6		P ₂	D ₂	D ₆	D ₁₀	D ₁₃
7		P ₃	D ₃	D ₇	P ₀	D ₀
8		D ₁	D ₂	D ₃	X	P ₀
9		D ₀	D ₁	X	P ₀	D ₀
10		D ₃	D ₆	D ₉	P ₁	D ₁
11		D ₄	D ₇	D ₁₀	P ₂	D ₂
12		D ₅	D ₈	.	.	.

Traditional Resilvering

- **Creating a new mirror (or RAID stripe):**
 - Copy one disk to the other (or XOR them together) so all copies are self-consistent – even though they're all random garbage!
- **Replacing a failed device:**
 - Whole-disk copy – even if the volume is nearly empty
 - No checksums or validity checks along the way
 - No assurance of progress until 100% complete – your root directory may be the last block copied
- **Recovering from a transient outage:**
 - Dirty region logging – slow, and easily defeated by random writes

Smokin' Mirrors

- **Top-down resilvering**
 - ZFS resilvers the storage pool's block tree from the root down
 - Most important blocks first
 - Every single block copy increases the amount of discoverable data
- **Only copy live blocks**
 - No time wasted copying free space
 - Zero time to initialize a new mirror or RAID-Z group
- **Dirty time logging (for transient outages)**
 - ZFS records the transaction group window that the device missed
 - To resilver, ZFS walks the tree and prunes by DTL
 - A five-second outage takes five seconds to repair

Ditto Blocks

- **Data replication above and beyond mirror/RAID-Z**
 - Each logical block can have up to three physical blocks
 - Different devices whenever possible
 - Different places on the same device otherwise (e.g. laptop drive)
 - All ZFS metadata 2+ copies
 - Settable on a per-file basis for precious user data (snv_61)
- **Detects and corrects silent data corruption**
 - If the first copy is missing or damaged, try the ditto blocks
 - In a multi-disk pool, ZFS survives any non-consecutive disk failures
 - In a single-disk pool, ZFS survives loss of up to 1/8 of the platter
- **ZFS survives failures that send other filesystems to tape**

Disk Scrubbing

- **Finds latent errors while they're still correctable**
 - ECC memory scrubbing for disks
- **Verifies the integrity of all data**
 - Traverses pool metadata to read every copy of every block
 - All mirror copies, all RAID-Z parity, and all ditto blocks
 - Verifies each copy against its 256-bit checksum
 - Self-healing as it goes
- **Minimally invasive**
 - Low I/O priority ensures that scrubbing doesn't get in the way
 - User-defined scrub rates coming soon
 - Gradually scrub the pool over the course of a month, a quarter, etc.

ZFS Scalability

- **Immense capacity (128-bit)**
 - Moore's Law: need 65th bit in 10-15 years
 - ZFS capacity: 256 quadrillion ZB (1ZB = 1 billion TB)
 - Exceeds quantum limit of Earth-based storage
 - Seth Lloyd, "Ultimate physical limits to computation."
Nature 406, 1047-1054 (2000)
- **100% dynamic metadata**
 - No limits on files, directory entries, etc.
 - No wacky knobs (e.g. inodes/cg)
- **Concurrent everything**
 - Byte-range locking: parallel read/write without violating POSIX
 - Parallel, constant-time directory operations

ZFS Performance

- **Copy-on-write design**
 - Turns random writes into sequential writes
 - Intrinsically hot-spot-free
- **Multiple block sizes**
 - Automatically chosen to match workload
- **Pipelined I/O**
 - Fully scoreboarded 24-stage pipeline with I/O dependency graphs
 - Maximum possible I/O parallelism
 - Priority, deadline scheduling, out-of-order issue, sorting, aggregation
- **Dynamic striping across all devices**
- **Intelligent prefetch**